

PATENT ABSTRACTS OF JAPAN

(11)Publication number : 2001-217720

(43)Date of publication of application : 10.08.2001

(51)Int.Cl.

H03M 7/30

(21)Application number : 2000-028359

(71)Applicant : INTERNATL BUSINESS MACH CORP
<IBM>

(22)Date of filing : 04.02.2000

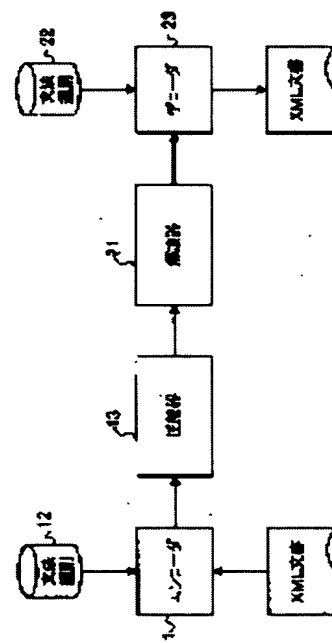
(72)Inventor : MARUYAMA HIROSHI
TAMURA TAKETO
URAMOTO NAOHIKO

(54) DATA COMPRESSING APPARATUS, DATA BASE SYSTEM, DATA COMMUNICATION SYSTEM, DATA COMPRESSING METHOD, STORAGE MEDIUM AND PROGRAM TRANSMITTER

(57)Abstract:

PROBLEM TO BE SOLVED: To provide a data compression for encoding the structural part of a document in a tree local language such as XML or ASN.1.

SOLUTION: In the data compressing apparatus for encoding and compressing data, there are provided a grammar rule 12 for the tree local language with which the data are expressed in a labeled tree structure, an encoder 11 for reading a document described in this tree local language, separating this document into a structure and contents and encoding this structure while using the grammar rule 12 and a compressor 13 for compression-encoding the contents of this document extracted by the encoder 11.



LEGAL STATUS

[Date of request for examination]

12.09.2000

[Date of sending the examiner's decision of rejection]

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number]

3368883

[Date of registration] 15.11.2002

[Number of appeal against examiner's decision of rejection]

[Date of requesting appeal against examiner's decision of rejection]

[Date of extinction of right]

Copyright (C); 1998,2003 Japan Patent Office

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号
特開2001-217720
(P2001-217720A)

(43) 公開日 平成13年8月10日 (2001.8.10)

(51) Int.Cl.⁷

H 0 3 M 7/30

識別記号

F I

H 0 3 M 7/30

テーマコード* (参考)

Z 5 J 0 6 4

審査請求 有 請求項の数13 O L (全 12 頁)

(21) 出願番号 特願2000-28359 (P2000-28359)

(22) 出願日 平成12年2月4日 (2000.2.4)

(71) 出願人 390009531

インターナショナル・ビジネス・マシー
ズ・コーポレーション

INTERNATIONAL BUSIN
ESS MACHINES CORPO
RATION

アメリカ合衆国10504、ニューヨーク州

アーモンク (番地なし)

(74) 代理人 100086243

弁理士 坂口 博 (外3名)

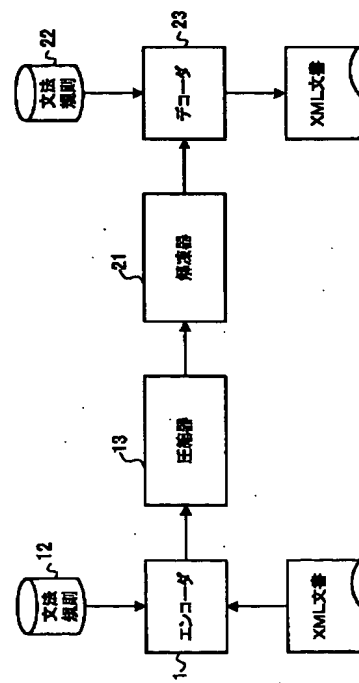
最終頁に続く

(54) 【発明の名称】 データ圧縮装置、データベースシステム、データ通信システム、データ圧縮方法、記憶媒体及びプログラム伝送装置

(57) 【要約】

【課題】 XMLやASN.1等の木ローカル言語における文書の構造部分の符号化を行うデータ圧縮を実現する。

【解決手段】 データを符号化して圧縮するデータ圧縮装置であって、データがラベル付きの木構造で表現される木ローカル言語の文法規則12と、この木ローカル言語で記述された文書を読み込んで、この文書を構造とコンテンツとに分け、文法規則12を用いてこの構造を符号化するエンコーダ11と、このエンコーダ11によって抽出されたこの文書のコンテンツを圧縮符号化する圧縮器13とを備える。



1

【特許請求の範囲】

【請求項 1】 データを符号化して圧縮するデータ圧縮装置であって、
データがラベル付きの木構造で表現される木ローカル言語の文法規則を格納した文法格納部と、
前記木ローカル言語で記述された文書を読み込んで、当該文書を構造とコンテンツとに分け、前記文法格納部に格納された前記文法規則を用いて当該構造を符号化するエンコーダと、
前記エンコーダによって抽出された前記文書のコンテンツを圧縮符号化する圧縮器とを備えることを特徴とするデータ圧縮装置。

【請求項 2】 前記エンコーダは、
処理対象である文書を構造とコンテンツとに分ける分割処理部と、
前記文法規則に対応するプッシュダウンオートマトンを構築するオートマトン構築部と、
前記オートマトン構築部により構築された前記プッシュダウンオートマトンを用いて、前記分割処理部により分割された前記文書の構造に対する構文解析を行い、当該構造の符号化されたデータ列を生成する符号化データ生成部とを備えることを特徴とする請求項 1 に記載のデータ圧縮装置。

【請求項 3】 前記エンコーダの前記符号化データ生成部は、前記オートマトン構築部により構築された前記プッシュダウンオートマトン中に存在する選択肢に対して符号を割り当て、当該プッシュダウンオートマトンを走らせて前記文書の構造を解析した際に当該選択肢の箇所で選択された選択肢に割り当てられている当該符号を出力することにより、当該構造の符号化データ列を生成することを特徴とする請求項 2 に記載のデータ圧縮装置。

【請求項 4】 前記圧縮器は、前記文書のコンテンツと共に、前記エンコーダにて符号化された当該文書の構造に対しても圧縮符号化を行うことを特徴とする請求項 1 に記載のデータ圧縮装置。

【請求項 5】 通信ネットワークを介してデータ送信を行う送信側データ処理装置と、当該送信側データ処理装置から送信されたデータを当該通信ネットワークを介して受信する受信側データ処理装置とを備えたデータ通信システムであって、
前記送信側データ処理装置は、
データがラベル付きの木構造で表現される木ローカル言語の文法規則を格納した第 1 の文法格納部と、
前記木ローカル言語で記述された送信文書を読み込んで、当該送信文書を構造とコンテンツとに分け、前記第 1 の文法格納部に格納された前記文法規則を用いて当該構造を符号化するエンコーダと、
前記エンコーダによって抽出された前記送信文書のコンテンツを圧縮符号化する圧縮器と、
前記エンコーダにより符号化された前記構造及び前記圧

2

縮器により圧縮符号化された前記コンテンツを送信する送信部とを備え、
前記受信側データ処理装置は、
前記送信側データ処理装置から送信された受信する受信部と、
前記送信側データ処理装置の前記第 1 の文法格納部に格納された文法規則と同一内容の文法規則を格納した第 2 の文法格納部と、
前記送信側データ処理装置の前記圧縮器による圧縮符号化手法に対応する解凍手法にて、前記受信部が受信した受信データのうち前記送信文書のコンテンツに対応するデータを解凍する解凍器と、
前記受信部が受信した受信データのうち前記送信文書の構造に対応するデータを、前記第 2 の文法格納部に格納された前記文法規則を用いて当該構造を復号化するデコーダとを備えることを特徴とするデータ通信システム。

【請求項 6】 データを記憶装置に格納して管理するデータベースシステムであって、
データがラベル付きの木構造で表現される木ローカル言語の文法規則を格納した文法格納部と、
前記木ローカル言語で記述された文書を読み込んで、当該文書を構造とコンテンツとに分け、前記文法格納部に格納された前記文法規則を用いて当該構造を符号化するエンコーダと、
前記エンコーダによって抽出された前記文書のコンテンツを圧縮符号化する圧縮器と、
前記エンコーダにより符号化された前記文書の構造と前記圧縮器により圧縮符号化された前記文書のコンテンツとを格納する記憶装置とを備えることを特徴とするデータベースシステム。

【請求項 7】 前記圧縮器は、前記文書のコンテンツと共に、前記エンコーダにて符号化された当該文書の構造に対しても圧縮符号化を行うことを特徴とする請求項 6 に記載のデータベースシステム。

【請求項 8】 データを符号化して圧縮するデータ圧縮方法であって、
データがラベル付きの木構造で表現される木ローカル言語で記述された文書を読み込んで、当該文書を構造とコンテンツとに分けるステップと、
前記木ローカル言語の文法規則を用いて、前記文書の構造を符号化するステップと、
前記文書のコンテンツを圧縮符号化するステップとを含むことを特徴とするデータ圧縮方法。

【請求項 9】 前記文書の構造を符号化するステップは、
前記文法規則に対応するプッシュダウンオートマトンを構築するステップと、
前記プッシュダウンオートマトン中に存在する選択肢に対して符号を割り当てるステップと、
前記プッシュダウンオートマトンを走らせて前記文書の

3

構造を深さ優先でたどりながら解析し、当該選択肢の箇所で選択された選択肢に割り当てられている前記符号を出力するステップと、

前記プッシュダウンオートマトンを走らせて出力された前記符号の列を前記文書の構造の符号化されたデータ列として出力するステップとを含むことを特徴とする請求項8に記載のデータ圧縮方法。

【請求項10】 前記文書の構造を符号化するステップに先だって、

処理対象である前記木ローカル言語の文書のノードに属性が含まれている場合に、当該属性を、当該属性を持つ要素の子ノードに変換することにより、前記木ローカル言語の文法規則及び前記文書を、前記プッシュダウンオートマトンで扱える木構造に変換するステップをさらに含むことを特徴とする請求項9に記載のデータ圧縮方法。

【請求項11】 前記文書の構造を符号化するステップの後に、

汎用的な他の圧縮符号化手法を用いて、当該符号化された当該文書の構造をさらに圧縮符号化するステップをさらに含むことを特徴とする請求項8に記載のデータ圧縮方法。

【請求項12】 コンピュータに実行させるプログラムを当該コンピュータの入力手段が読取可能に記憶した記憶媒体において、

前記プログラムは、

データがラベル付きの木構造で表現される木ローカル言語で記述された文書を読み込んで、当該文書を構造とコンテンツとに分ける処理と、

前記木ローカル言語の文法規則を用いて、前記文書の構造を符号化する処理と、

前記文書のコンテンツを圧縮符号化する処理とを前記コンピュータに実行させることを特徴とする記憶媒体。

【請求項13】 コンピュータに、

データがラベル付きの木構造で表現される木ローカル言語で記述された文書を読み込んで、当該文書を構造とコンテンツとに分ける処理と、前記木ローカル言語の文法規則を用いて、前記文書の構造を符号化する処理と、前記文書のコンテンツを圧縮符号化する処理とを実行させるプログラムを記憶する記憶手段と、

前記記憶手段から前記プログラムを読み出して当該プログラムを送信する送信手段とを備えたことを特徴とするプログラム伝送装置。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、XMLやASN.1等の木ローカル言語で記述された文書を圧縮するデータ圧縮方法に関する。

【0002】

【従来の技術】XML (Extensible Markup Language)

4

は、文書の意味構造を簡単なマークで記述(マークアップ)するマークアップ言語の一種である。XMLでは、文法を規定し、文書の構成要素に論理的意味を持たせることにより、ユーザが独自の拡張を行うことが可能である。そのため、インターネットにおけるデータ交換に用いるデータフォーマットとして期待されている。

【0003】XMLには、DTD (文書型定義: Document Type Definition) という概念があり、あるDTDに関して、文書が妥当 (valid) かどうかを判定できる。

10 具体的には、例えば、<BOOK>というノードの下には、<TITLE>、<AUTHOR>、<PUBLISHER>というノードがこの順にそれぞれ1回ずつ現れる、というような文法規則を規定する。そして、所定の文書が妥当かどうか、すなわち、当該文法規則に則しているかどうかを判定することができる。

【0004】XML文書の構造は、形式言語的には、木言語のうちの木ローカルと呼ばれるクラスに属する。木ローカル言語とは、データがラベル付きの木構造で表現され、かつ、正しいデータが木の各ノードラベルに対して子ノードのラベルの正規言語で規定されるような木言語である。すなわち、木ローカル言語では、所定の文法 (XMLにおいてはDTDにて規定される) に属する木の集合が、各ノードの子ノードのリストを規定する正規言語によって決定付けられる。この種の木ローカル言語としては、他に、ASN.1 (Abstract Syntax Notation 1) 等がある。

【0005】近年、XMLを用いて、ネットワーク上のリソースやHTMLでは実現不可能であった複雑なデータ構造を記述し、ビジネスアプリケーションを構築する動きが広がっている。このようなシステムでは、大規模なXML文書がアプリケーション間で交換されることが予想される。一般に、データ交換を行ったりデータベースにデータを格納したりする際には、送信効率の向上や格納するデータファイルのサイズを縮小するために、データファイルの圧縮を行う。そのため、種々のデータフォーマットに対して適用できる汎用的なデータ圧縮技術や、特定のデータフォーマットに特化したデータ圧縮技術などが数多く提案されている。したがって、XML文書の交換においてもこれらのデータ圧縮技術を適用することが考えられる。なお、XML文書は、タグ名などかなりの冗長度があるため、高い圧縮率が期待できるデータ構造であるといえる。

【0006】

【発明が解決しようとする課題】上述したように、データ交換を行ったりデータベースにデータを格納したりする際には、データファイルの圧縮を行うことが一般的である。そして、XML等の木ローカル言語では、タグなどの文書構造を示す情報部分に対して高い圧縮率を期待できる。

50 【0007】例えば、データ通信において、通信の当事

5

者が文法Gを共有しており、互いに当該文法Gに対して妥当なXML文書のみをやりとりすることが保証される場合を考える。また、上記のように文法Gにおいて、〈BOOK〉というノードの下には、〈TITLE〉、〈AUTHOR〉、〈PUBLISHER〉というノードがこの順にそれぞれ1回ずつ現れる、という文法規則が規定されているものとする。この場合、XML文書の受信者は、受信したXML文書の中で、例えば〈BOOK〉というタグを検出した場合に、その最初の子ノードは〈TITLE〉であると予測することができる。したがって、〈TITLE〉という情報は、この仮定のもとでは冗長である。この考えを用いることで、メッセージ長を最適化するような、XML文書のエンコーディング方法が可能である。このことは、XMLに限らず、任意の木ローカル言語（例えば、ASN.1）に対しても同様である。しかしながら、従来、このような木ローカル言語における文書の構造部分を符号化する圧縮方法は何ら提案されていなかった。

【0008】そこで本発明は、XMLやASN.1等の木ローカル言語における文書の構造部分の符号化を行うデータ圧縮を実現する。

【0009】また本発明は、木ローカル言語に特化すると共に、他の汎用的なデータ圧縮技術と併用することにより、高い圧縮率を期待できるデータ圧縮方法を提供する。

【0010】

【課題を解決するための手段】かかる目的のもと、本発明は、データを符号化して圧縮するデータ圧縮装置であって、データがラベル付きの木構造で表現される木ローカル言語の文法規則を格納した文法格納部と、この木ローカル言語で記述された文書を読み込んで、この文書を構造とコンテンツとに分け、文法格納部に格納された文法規則を用いてこの構造を符号化するエンコーダと、このエンコーダによって抽出されたこの文書のコンテンツを圧縮符号化する圧縮器とを備えることを特徴としている。木ローカル言語とは、データがラベル付きの木構造で表現され、かつ、正しいデータが木の各ノードラベルに対して子ノードのラベルの正規言語で規定されるような木言語である。

【0011】ここで、このエンコーダは、処理対象である文書を構造とコンテンツとに分ける分割処理部と、文法規則に対応するプッシュダウンオートマトンを構築するオートマトン構築部と、このオートマトン構築部により構築された前記プッシュダウンオートマトンを用いて、前記分割処理部により分割された前記文書の構造に対する構文解析を行い、当該構造の符号化されたデータ列を生成する符号化データ生成部とを備えることを特徴としている。

【0012】さらに、このエンコーダの符号化データ生成部は、オートマトン構築部により構築されたこのプッシュダウンオートマトン中に存在する選択肢に対して符

6

号を割り当て、このプッシュダウンオートマトンを走らせて木ローカル言語で記述された文書の構造を解析した際にこの選択肢の箇所を選択された選択肢に割り当てられている符号を出力することにより、この構造の符号化データ列を生成することを特徴としている。このような構成とすれば、タグなどのラベルを用いて記述されたこの文書の構造を簡単な符号列に変換（符号化）することが可能となる。プッシュダウンオートマトンを用いて文書の構造を解析する際には、この文書の木構造を深さ優先でたどりながら解析を行う。すなわち親ノードからの距離が等しい階層ごとに解析していくのではなく、深さ方向のノードどうしのつながり（親子関係）を優先させて木をたどりながら解析する。

【0013】またここで、圧縮器は、この木ローカル言語で記述された文書のコンテンツと共に、エンコーダにて符号化されたこの文書の構造に対しても圧縮符号化を行うことを特徴としている。この圧縮器による圧縮手法は、特に制限はなく、従来から用いられている汎用的な手法を採用することができる。エンコーダにより文書の構造を符号化した結果、ある程度規則的なデータ列が得られる場合は、この符号化データ列に対し、PKZIP等の汎用的な手法を用いてさらに圧縮符号化を行うことにより、高い圧縮率を期待できる。そこで、文書のコンテンツを圧縮する際に、この符号化データ列を合わせて圧縮することが好ましい。さらに、同じ文書における構造の符号化データ列とコンテンツとを結合した上で圧縮を行うことにより、構造とコンテンツとが別ファイルとなることを避けられるので、ファイル管理の上でも好ましい。

【0014】また、本発明は、通信ネットワークを介してデータ送信を行う送信側データ処理装置と、この送信側データ処理装置から送信されたデータをこの通信ネットワークを介して受信する受信側データ処理装置とを備えたデータ通信システムであって、この送信側データ処理装置は、データがラベル付きの木構造で表現される木ローカル言語の文法規則を格納した第1の文法格納部と、この木ローカル言語で記述された送信文書を読み込んで、この送信文書を構造とコンテンツとに分け、第1の文法格納部に格納された文法規則を用いてこの送信文書の構造を符号化するエンコーダと、このエンコーダによって抽出されたこの送信文書のコンテンツを圧縮符号化する圧縮器と、エンコーダにより符号化された構造及び圧縮器により圧縮符号化されたコンテンツを送信する送信部とを備え、この受信側データ処理装置は、この送信側データ処理装置から送信された受信する受信部と、送信側データ処理装置の第1の文法格納部に格納された文法規則と同一内容の文法規則を格納した第2の文法格納部と、送信側データ処理装置の圧縮器による圧縮符号化手法に対応する解凍手法にて、この受信部が受信した受信データのうち送信文書のコンテンツに対応するデー

データを解凍する解凍器と、この受信部が受信した受信データのうち送信文書の構造に対応するデータを、第2の文法格納部に格納された前記文法規則を用いて当該構造を復号化するデコーダとを備えることを特徴としている。このように、データの送信側と受信側とで予め共通の文法規則を用意しておけば、データ通信において、木ローカル言語で記述された文書に対するきわめて圧縮率の高い圧縮を行うことができ、通信効率を向上させることができる点できわめて優れている。なお、ビジネス間通信においては、木ローカル言語の文法規則として共通のものを用いることが予め定められることが一般的なので、本発明を導入することは容易である。

【0015】さらにまた、本発明は、データを記憶装置に格納して管理するデータベースシステムであって、データがラベル付きの木構造で表現される木ローカル言語の文法規則を格納した文法格納部と、この木ローカル言語で記述された文書を読み込んで、この文書を構造とコンテンツとに分け、文法格納部に格納された文法規則を用いてこの文書の構造を符号化するエンコーダと、このエンコーダによって抽出された文書のコンテンツを圧縮符号化する圧縮器と、エンコーダにより符号化されたこの文書の構造と圧縮器により圧縮符号化されたこの文書のコンテンツとを格納する記憶装置とを備えることを特徴としている。

【0016】ここで、この圧縮器は、前記文書のコンテンツと共に、前記エンコーダにて符号化された当該文書の構造に対しても圧縮符号化を行うことを特徴としている。同じ文書における構造の符号化データ列とコンテンツとを結合した上で圧縮を行うことにより、高い圧縮率を期待でき、さらに構造とコンテンツとが別ファイルとなることを避けられるので、ファイル管理の上でも好ましい。

【0017】また、本発明は、データを符号化して圧縮するデータ圧縮方法であって、データがラベル付きの木構造で表現される木ローカル言語で記述された文書を読み込んで、この文書を構造とコンテンツとに分けるステップと、この木ローカル言語の文法規則を用いて、この文書の構造を符号化するステップと、この文書のコンテンツを圧縮符号化するステップとを含むことを特徴としている。

【0018】ここで、文書の構造を符号化するステップは、文法規則に対応するプッシュダウンオートマトンを構築するステップと、プッシュダウンオートマトン中に存在する選択肢に対して符号を割り当てるステップと、このプッシュダウンオートマトンを走らせて前記文書の構造を深さ優先でたどりながら解析し、この選択肢の箇所で選択された選択肢に割り当てられているこの符号を出力するステップと、このプッシュダウンオートマトンを走らせて出力されたこの符号の列をこの木ローカル言語で記述された文書の構造の符号化されたデータ列とし

て出力するステップとを含むことを特徴としている。このような構成とすれば、タグなどのラベルを用いて記述されたこの文書の構造を簡単な符号列に変換（符号化）することが可能となる。

【0019】このデータ圧縮方法において、木ローカル言語で記述された文書の構造を符号化するステップに先だって、処理対象であるこの木ローカル言語の文書のノードに属性が含まれている場合に、この属性を、この属性を持つ要素の子ノードに変換することにより、この木ローカル言語の文法規則及び文書を、プッシュダウンオートマトンで扱える木構造に変換するステップをさらに含むことを特徴としている。このように構成すれば、XMLのように処理対象の文書に属性が含まれている場合にも、プッシュダウンオートマトンを用いた構造の符号化を行うことができる点で好ましい。

【0020】さらに、文書の構造を符号化するステップの後に、汎用的な他の圧縮符号化手法を用いて、符号化された文書の構造をさらに圧縮符号化するステップをさらに含むことを特徴としている。このような構成とすることにより、さらに高い圧縮率を期待できる点で好ましい。

【0021】また、本発明は、コンピュータに実行させるプログラムを当該コンピュータの入力手段が読取可能に記憶した記憶媒体において、このプログラムは、データがラベル付きの木構造で表現される木ローカル言語で記述された文書を読み込んで、この文書を構造とコンテンツとに分ける処理と、この木ローカル言語の文法規則を用いて、この文書の構造を符号化する処理と、この文書のコンテンツを圧縮符号化する処理とをこのコンピュータに実行させることを特徴としている。このような構成とすれば、このプログラムをインストールしたあらゆる情報処理装置において、この木ローカル言語にて記述された文書を高い圧縮率で圧縮することができ、通信効率や記憶効率を向上させることができる。

【0022】さらにまた、本発明は、コンピュータに、データがラベル付きの木構造で表現される木ローカル言語で記述された文書を読み込んで、この文書を構造とコンテンツとに分ける処理と、この木ローカル言語の文法規則を用いて、この文書の構造を符号化する処理と、この文書のコンテンツを圧縮符号化する処理とを実行させるプログラムを記憶する記憶手段と、この記憶手段からこのプログラムを読み出してこのプログラムを送信する送信手段とを備えたことを特徴としている。このような構成とすれば、このプログラム伝送装置からこのプログラムをダウンロードしてインストールしたあらゆる情報処理装置において、この木ローカル言語にて記述された文書を高い圧縮率で圧縮することができ、通信効率や記憶効率を向上させることができる。

【0023】

【発明の実施の形態】以下、添付図面に示す実施の形態

に基づいてこの発明を詳細に説明する。図1は、本実施の形態における文書圧縮システムの全体構成を説明する図である。図1において、符号11はエンコーダであり、圧縮対象である文書を構造とコンテンツとに分け、その構造部分を、所定の記憶装置に記憶された文法規則12を用いて符号化する。符号13は圧縮器であり、エンコーダ11により符号化された構造部分と、当該文書のコンテンツ部分とを圧縮する。符号21は解凍器であり、圧縮器13により圧縮された文書を解凍する。解凍器21により解凍された時点では、当該文書は、コンテンツ部分とエンコーダ11により符号化された構造部分とに分かれている。符号23はデコーダであり、符号化されている構造部分を、所定の記憶装置に記憶された文法規則22を用いて複合化し、コンテンツ部分と合わせて、文書を復元する。

【0024】本実施の形態をデータ通信の際のデータ圧縮に用いる場合は、エンコーダ11及び圧縮器13は送信側に置かれ、解凍器21及びデコーダ23は受信側に置かれることとなる。また、データベースシステムにおいて格納するデータファイルを圧縮するために用いる場合は、データの流れに応じて、エンコーダ11がデコーダ23として動作し、圧縮器13が解凍器21として動作する。

【0025】以下、処理対象の木ローカル言語としてXMLを用いる場合を例として説明する。図2は、本実施の形態によるデータ圧縮の手順を説明する図である。図2に示す本実施の形態によるデータ圧縮処理では、まず、処理対象であるXML文書201が、エンコーダ11に読み込まれて、構造202とコンテンツ204とに分解される。ここで、構造202とは、当該XML文書201の木構造、タグ名及び属性名であり、コンテンツ204とは、当該XML文書201の#PCDATA及び属性値である。XML文書201を構造202とコンテンツ204とに分解するのは、構造202とコンテンツ204とは、通常全く異なる統計的偏りを持っているため、独立して圧縮した方が効果的だからである。

【0026】次に、XML文書201を分解して得られた構造202が、エンコーダ11により符号化される。構造202の符号化には、文法規則12が用いられる。ここでは、処理対象をXML文書としているので、文法規則12はDTDにより規定される。符号化処理の詳細な内容については後述する。符号化された結果である符号化データ列203とコンテンツ204とは、圧縮器13に送られる。

【0027】最後に、圧縮器13において、符号化データ列203とコンテンツ204とが圧縮符号化され、これらのデータを合わせて、圧縮されたXML文書205が生成される。圧縮器13による符号化処理には、LZ77等の既知の適当な手法を用いる。ここで、圧縮器13による圧縮符号化は、主としてコンテンツ204に対

して行われることとなる。しかしながら、PKZIP等の汎用的な圧縮符号化方式は、符号化データ列203に対しても有効である。後述するように、本実施例では、符号化データ列203は数字列として出力される。したがって、ある程度規則的な数字列であるような場合は、高い圧縮率が期待できる。そこで、圧縮器13では、コンテンツ204と共に符号化データ列203も圧縮符号化を行う。なお、符号化データ列203に対して圧縮器13による圧縮を行うか否かは任意である。符号化データ列203と圧縮器13により圧縮されたコンテンツ204とを単に関連づけ、または結合して、送受信したり、記憶装置に格納したりするようにしても良い。以上のように、本実施の形態は、XML文書201における構造202の部分を本実施の形態による手法を用いて圧縮し、さらに符号化された構造202の部分とコンテンツ204の部分とを、従来の手法を用いて圧縮する。このため、本実施の形態によるデータ圧縮は、全体としては、従来の種々の圧縮手法と併用して行うこととなる。

【0028】以上のようにして圧縮されたXML文書205を解凍する場合は、上記の圧縮課程の反対の過程をたどる。すなわち、まず解凍器21において、圧縮器13が圧縮符号化に用いた手法に対応する手法で、符号化データ列203とコンテンツ204とを解凍する。次に、デコーダ23において、解凍された符号化データ列203を、文法規則12と同一の文法規則22を用いて複合化する。復号化処理の詳細な内容については後述する。文法規則22はDTDにより規定される。そして、この復号化処理により得られた構造202と解凍器21により解凍されたコンテンツ204とを用いて、XML文書201を復元する。

【0029】次に、本実施の形態によるXML文書の構造に対する符号化処理の内容を詳細に説明する。ここでは、処理対象のXML文書は、簡単のため、属性を含まないものとし、XML文書中のすべての実体は展開されているものとする。属性の扱いについては後述する。

【0030】図3は、XML文書の構造を符号化するエンコーダ11の構成を説明する機能ブロック図である。図3を参照すると、エンコーダ11は、処理対象であるXML文書201を構造202とコンテンツ204とに分割する分割処理部111と、文法規則12に基づいて後述するプッシュダウンオートマトンを構築するオートマトン構築部112と、オートマトン構築部112により構築されたプッシュダウンオートマトンを符号化トランスデューサとして用いて構造202の符号化データ列203を生成する符号化データ列生成部113とを備える。

【0031】図4に処理対象とするXML文書の例を示す。XML文書のコンテンツとは、内容モデルにおける、#PCDATAに相当する部分の文字列のリストである。すなわち、図4のXML文書におけるコンテンツ

10

20

30

40

50

は、“String1”、“String2”、“String3”、“String4”という4つの文字列からなるリストである。これは、例えば、ナル文字で終わる文字列を並べたバイト列として、以下のように、コンパクトに表現可能である(但し“¥0”はナル文字をあらわす)。“String1¥0String2¥0String3¥0String4¥0”この文字列は、上述したように、構造部分とは別に圧縮符号化される。また、図4のXML文書における構造は、図5に示すようになる。これは、図4のXML文書中のコンテンツに相当する文字列をプレースホルダ(□)に置き換えたものである。

【0032】本実施の形態において、エンコーダ11は、分割処理部111により図4に示すXML文書から図5に示す構造を取り出し、オートマトン構築部112により文法規則12を用いてプッシュダウンオートマトンを構築し、符号化データ列生成部113により当該プッシュダウンオートマトンを用いて当該構造を符号化する。図6は、文法規則12を規定するDTDの例を示す。分割処理部111による分割処理の後、オートマトン構築部112は、文法規則12を用いた符号化処理のために、当該DTDに対応するプッシュダウンオートマトンを構築する。図6のDTDによれば、要素Aが現れた場合は、次に、要素B、要素Cがこの順番で1回ずつ現れた後に終了することを示す。同様に、要素Bが現れたときは、次に、要素Dが1回現れた後に終了することを示す。また、要素Cが現れたときは、次に、要素Eまたは要素Fが0回以上現れた後に終了することを示す。さらにまた、要素Eが現れたときは、次に、要素Gまたは要素Hのいずれか一方が1回現れた後に終了することを示す。図7は、図6に示したDTDに対応する自然なプッシュダウンオートマトンを示す図である。なお、非終端記号DとGに関しては、終端記号#PCDATAを取るだけの自明なルールなので、省略してある。このようなオートマトンは、文法の各非終端記号について曖昧さ無く構築できる。したがって、本実施の形態をデータ通信に用いる場合、送信側と受信側との共通のDTDからは、全く同一のプッシュダウンオートマトンを構築することができる。

【0033】通常、プッシュダウンオートマトンは、入力列の構文解析を行うために用いられる。その意味では、このプッシュダウンオートマトンは、表層のシンボル列、すなわち、#PCDATA(またはプレースホルダ“□”)の1個以上の並びからなる全ての列を受理する。しかし、解析の結果生成される構文解析木としては、例えば、構文木のノードAの子供として、ノードB、ノードCがこの順に現われなければならない。また、要素Cの後には、空遷移で最終状態に遷移する。このように、このプッシュダウンオートマトンは、解析済みのXML文書(例えば、DOM木)のような、構文解析木が文法を満たすかどうかのチェックに用いることもできる。

【0034】図8に示す構文木を例として、プッシュダウンオートマトンを用いた文法のチェックについて説明する。なお、図8において、各リーフの#PCDATAは省略してある。この構文木が図6のDTDに規定される文法によって生成可能かどうかを調べるには、この構文木の各ノードに対して、そのノードの非終端記号に対応するオートマトンによって、その子ノードの列が受理できるかどうかを調べればよい。例えば、要素Aは、子ノードとしてBCという列を持っている。これは、非終端記号Aに対応するオートマトンによって受理される

(図7のA参照)。したがって、この部分については文法を満たすことがわかる。同様に、全てのノードについて、対応するオートマトンを使ってPreorderでトラバース(深さ優先でたどる)すれば、文法のチェックは終了する。構文解析木に対して、プッシュダウンオートマトンをこのように使うことを、以下の説明では妥当性検証と呼ぶ。なお、以上の操作で用いる各非終端記号に対するオートマトンについては、終了状態へのε遷移を除き、決定性で最小のオートマトンであるものとする。

【0035】次に、オートマトン構築部112は、図7のプッシュダウンオートマトンを、XML文書の構造部分(図5参照)を符号化するトランスデューサ、すなわち、文字列の構文解析のオートマトンに変換する。図7のプッシュダウンオートマトンにおいて、入力を、4つの#PCDATA(またはプレースホルダ“□”)からなる並びとし、開始記号をAとして解析を開始すると、ノードA、ノードB、ノードDが順に作られ、最初の#PCDATAが認識される。次に、ノードCが作られたところで選択肢が発生する。すなわち、ノードEを作るべきか、ノードFを作るべきか、それともノードCをこれで終わりにして上位のノードに戻るか、の3通りである。そこで、これらの3通りの選択肢に、ラベルのアルファベット順に1、2、3という番号を割り振る

(ラベルεは常に最後と決める)。同様に、ルールEの最初の状態も、ノードGを作るかまたはノードHを作るという選択肢を持つので、これらについても、1、2という番号を割り当てる。なお、ここでは選択肢に番号を割り振るとしたが、選択肢の識別に用いる符号は数字に限らない。アルファベットや記号など任意の符号を用いて選択肢を特定することが可能である。図9は、図7のプッシュダウンオートマトンから変換された符号化トランスデューサである。

【0036】エンコーダ11の符号化データ列生成部113は、オートマトン構築部112により構築された符号化トランスデューサを走らせる。図9に示す符号化トランスデューサは、妥当性検証(Preorderでのルール適用)を行った場合に、対応する選択肢番号があれば、その番号を出力する。すなわち、図9において、ルールA、B、F、Hに関しては、選択肢がないので何ら出力を行わず、ルールC及びルールDが適用された場合は、

該当する番号を出力する。例えば、図 8 の構文木に対して妥当性検証を行った場合、木をたどるにしたがって、図 10 に示すような番号を出力する。以上の処理により、“112123”という番号列が得られる。この番号列は、プッシュダウンオートマトンの動きを厳密に規定している。したがって、この番号列を、図 4 に示した XML 文書の構造部分 (図 5) の符号化されたデータ列として扱うことができる。

【0037】次に、本実施の形態による XML 文書の構造に対する復号化処理の内容を説明する。上記の手順を経て符号化された XML 文書を復号するには、符号化トランスデューサの入出力を逆にして適用すればよい。したがって、デコーダ 23 は、図 7 と同一のプッシュダウンオートマトンを用いて復号化トランスデューサを生成し、復号化処理を実行する。上述したように、オートマトンは、文法の各非終端記号について曖昧さ無く構築できるため、DTD にて規定される文法規則 12 と文法規則 22 とが共通であれば、デコーダ 23 においても図 7 と全く同一のプッシュダウンオートマトンを構築することができる。

【0038】図 11 は、図 7 と同一のプッシュダウンオートマトンから変換された復号化トランスデューサである。図 11 に示す復号化トランスデューサにおいて、i/B という表現は、「i という入力文字列を見たら、B というルールを呼び、その後、次の状態へ遷移する」という遷移をあらわす。これにより、エンコーダ 11 から出力された番号列を入力して、対応する構文解析木を生成する。上述した“112123”という番号列を入力した場合、元々の番号の割り当て方から、このプッシュダウンオートマトン (復号化トランスデューサ) は、曖昧さ無く XML 文書の符号化番号列を受理することができる。したがって、生成される構文解析木は、図 8 に示したオリジナルの構文解析木と同一となる。これにより、当該 XML 文書の構造部分が復元されることとなる。

【0039】次に、属性の扱いについて説明する。本実施の形態では、属性は、プッシュダウンオートマトンで扱えるように、木構造に変換する。具体的には、属性を持つ全ての要素 (ELEMENT) に関して、それらの属性を子ノードとして取るように変換する。この際、属性は属性名のアルファベット順に現れるものとする。そして、#REQUIRED である属性は、そのまま並べる。また、#IMPLIED である属性は、オプションである“?”をつける。なお、#FIXED である属性については、元々情報がないので、変換後の DTD には含めない。図 12 は、例示的に、所定の DTD における変換前後の状態を比較する図である。図 12 において、左側に示される DTD は、右側に示される形式に変換される。このような DTD に応じて、属性を含む XML 文書においても、属性を要素に変換する。図 13 は、例示的

に、所定の XML 文書における変換前後の状態を比較する図である。

【0040】以上のようにして、DTD と XML 文書とを、属性を持たない状態とした後に、上述した符号化処理及び復号化処理を実行する。なお、DTD の変換は、プッシュダウンオートマトンを構築する前の段階で予め行っても良いし、プッシュダウンオートマトンを構築する段階で逐次行っても良い。前者の場合は、変換後の新しい DTD を用いてプッシュダウンオートマトンを構築する手順となる。また、後者の場合は、元の (属性を持つ) DTD から直接プッシュダウンオートマトンを構成する手順となる。

【0041】以上説明したように、本実施の形態は、XML 文書を圧縮する側と解凍する側の双方において同一の DTD を共有することが不可欠である。したがって、本実施の形態をデータ通信の際のデータ圧縮に用いる場合は、送信側データ処理装置と受信側データ処理装置とにそれぞれ同一の DTD を用意する必要がある。図 14 は、本実施の形態を用いたデータ通信システムの構成例を説明する図である。送信側データ処理装置 1410 において、エンコーダ 11 は、データ処理部から XML 文書を受け取り、DTD 1411 (図 1 の文法規則 12 に相当) を用いて構造部分の符号化を行う。圧縮器 13 は、符号化された構造部分及びコンテンツ部分の圧縮を行う。送信部 1412 は、エンコーダ 11 及び圧縮器 13 にて圧縮された XML 文書を、通信ネットワークを介して受信側データ処理装置 1420 へ送信する。受信側データ処理装置 1420 において、受信部 1422 は通信ネットワークを介して受信した受信データを解凍器 21 へ送る。解凍器 21 は、受け取った受信データの解凍を行う。この時点で XML 文書のコンテンツ部分は復元される。デコーダ 23 は、解凍された受信データの構造部分の符号化データ列を、DTD 1421 (図 1 の文法規則 22 に相当) を用いて復号化する。そして、解凍されているコンテンツ部分と合わせて XML 文書を復元し、データ処理部へ送る。ここで、送信側データ処理装置 1410 の DTD 1411 と受信側データ処理装置 1420 の DTD 1421 とが同一の内容となっている。なお、電子商取引等のビジネス間通信では、アプリケーションどうして XML 文書をやりとりする場合、DTD は事前に合意されている場合がほとんどである。したがって、DTD を共有していることを前提に、本実施の形態をビジネス間通信に利用することができる。

【0042】また、データベースシステムにおいて格納するデータファイルを圧縮するために用いる場合は、XML ファイルの構造を符号化するために用いた DTD を、そのまま復号化する際に利用できるため、DTD が共有されているかどうかを考慮する必要はない。図 15 は、本実施の形態を用いたデータベースシステムの構成を説明する図である。データベースシステム 1500 に

において、エンコーダ 11 は、データ処理部から XML 文書を受け取り、DTD 1501 (図 1 の文法規則 12 に相当) を用いて構造部分の符号化を行う。圧縮器 13 は、符号化された構造部分及びコンテンツ部分の圧縮を行う。そして、エンコーダ 11 及び圧縮器 13 にて圧縮された XML 文書が記憶装置 1502 に格納される。記憶装置 1502 に格納されている圧縮された XML 文書を読み出す場合は、圧縮器 13 が解凍器 21 として動作し、エンコーダ 11 がデコーダ 23 として動作する。XML 文書の構造部分の復号化には、符号化の際に用いた DTD 1501 を再度使用する。

【0043】なお、以上の説明では、木ローカル言語として XML を用いた場合を例として説明したが、ASN.1 等の他の木ローカル言語においてもそのまま利用することができる。ただし、この場合においても、データファイルを圧縮する側と解凍する側とで、上述した XML における DTD のような文法規則を共有することが必要である。

【0044】

【発明の効果】以上説明したように、本発明によれば、木ローカル言語における文書の構造部分の符号化を行うデータ圧縮を実現することができる。

【0045】また、木ローカル言語に特化すると共に、他の汎用的なデータ圧縮技術と併用することにより、高い圧縮率を期待できるデータ圧縮方法を提供することができる。

【図面の簡単な説明】

【図 1】 本実施の形態における文書圧縮システムの全体構成を説明する図である。

【図 2】 本実施の形態によるデータ圧縮の手順を説明する図である。

【図 3】 本実施の形態におけるエンコーダの構成を説

明する図である。

【図 4】 本実施の形態の処理対象である XML 文書の例を示す図である。

【図 5】 図 4 の XML 文書における構造を示す図である。

【図 6】 本実施の形態にて用いる文法規則の例を示す図である。

【図 7】 図 6 の文法規則から構築されるプッシュダウンオートマトンを示す図である。

【図 8】 プッシュダウンオートマトンを用いた文法のチェックの手法を説明するための構文木を例示的に示す図である。

【図 9】 図 7 のプッシュダウンオートマトンを用いて生成された符号化トランスデューサを示す図である。

【図 10】 図 8 の構文木に対して妥当性検証を行った結果を例示的に示す図である。

【図 11】 図 7 と同一のプッシュダウンオートマトンを用いて生成された復号化トランスデューサを示す図である。

【図 12】 属性を持つ DTD を、属性を持たない DTD に変換した状態を説明する図である。

【図 13】 属性を持つ XML 文書を、属性を持たない XML 文書に変換した状態を説明する図である。

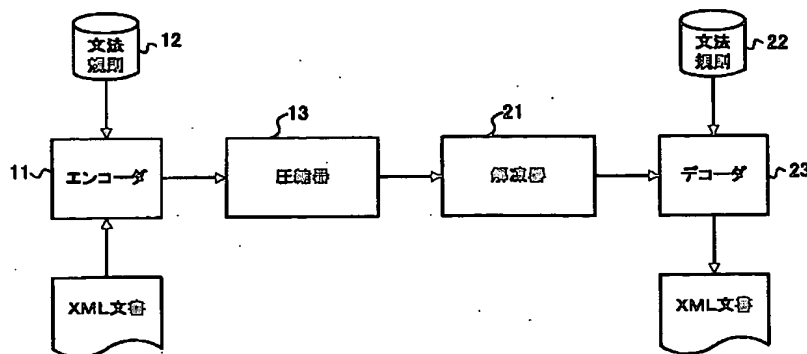
【図 14】 本実施の形態をデータ通信システムに適用した場合の構成を説明する図である。

【図 15】 本実施の形態をデータベースシステムに適用した場合の構成を説明する図である。

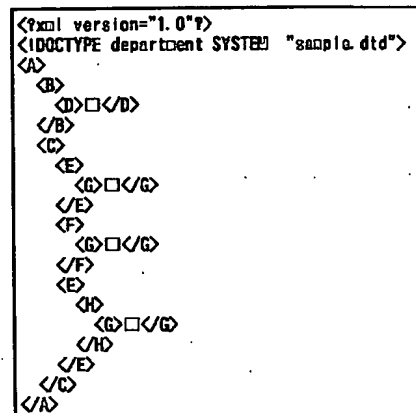
【符号の説明】

11…エンコーダ、12…文法規則、13…圧縮器、21…解凍器、23…デコーダ、201…XML 文書、202…構造、203…符号化データ列、204…コンテンツ、205…圧縮された XML 文書

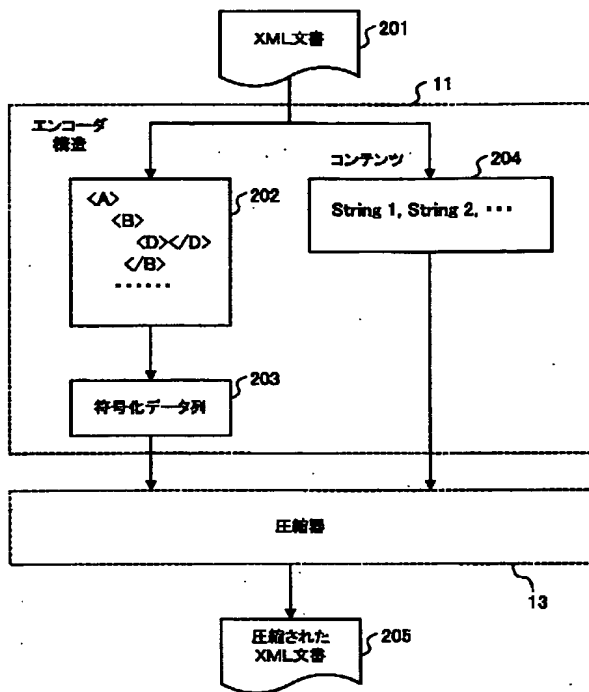
【図 1】



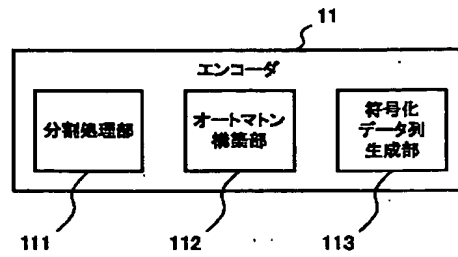
【図 5】



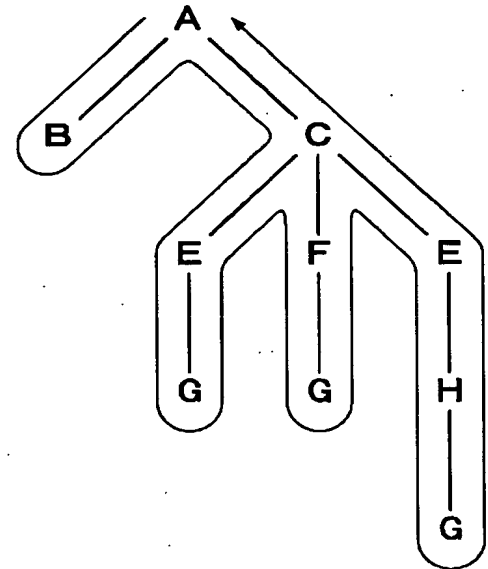
【図 2】



【図 3】



【図 8】



【図 4】

```

<?xml version="1.0"?>
<!DOCTYPE department SYSTEM "sample.dtd">
<A>
  <B>
    <D>String1</D>
  </B>
  <C>
    <E>String2</E>
  </C>
  <F>
    <G>String3</G>
  </F>
  <E>
    <H>
      <G>String4</G>
    </H>
  </E>
</A>

```

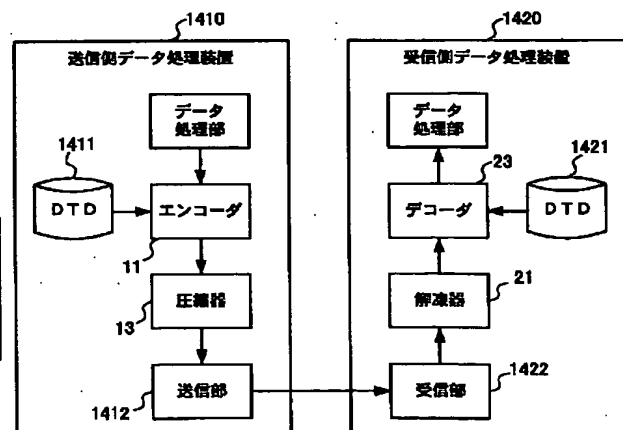
【図 6】

```

<ELEMENT A (B, C)>
<ELEMENT B (D)>
<ELEMENT C (E|F)*>
<ELEMENT E (G|H)>
<ELEMENT F (G)>
<ELEMENT H (G)>
<ELEMENT D (#PCDATA)>
<ELEMENT G (#PCDATA)>

```

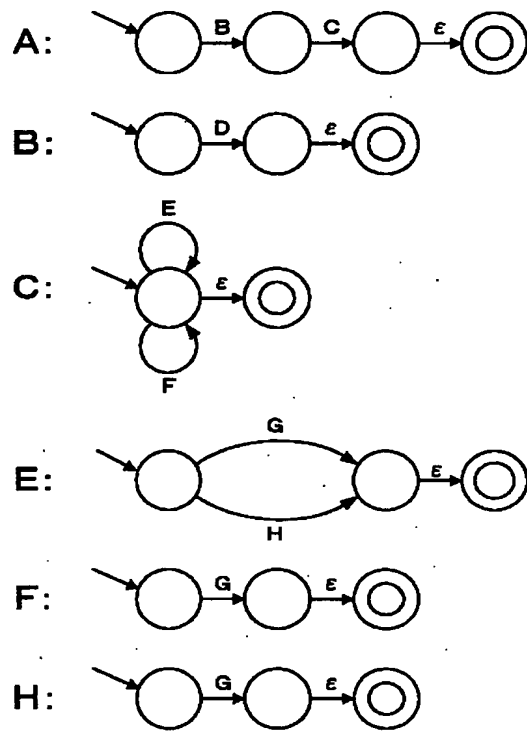
【図 14】



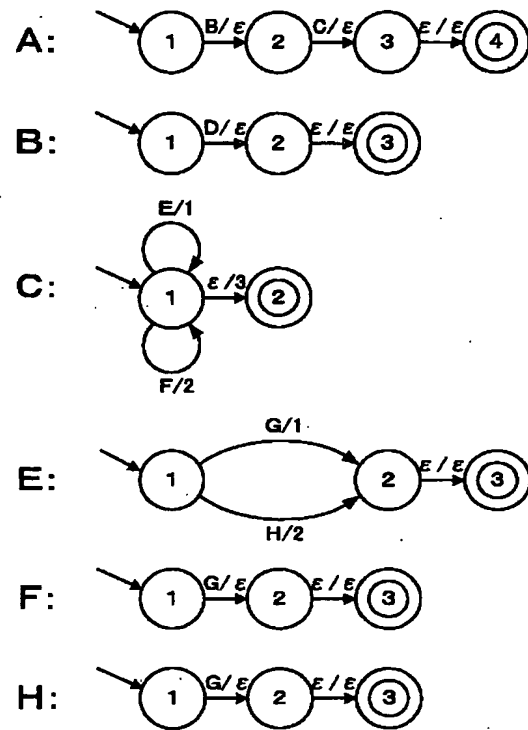
【図 13】

変換前	変換後
<pre> <department serialNo="012345" manager="yes"> <name>aaaa</name> <email>mail@address.com</email> </department> </pre>	<pre> <department <serialNo>012345</serialNo> <manager>yes</manager> <name>aaaa</name> <email>mail@address.com</email> </department> </pre>

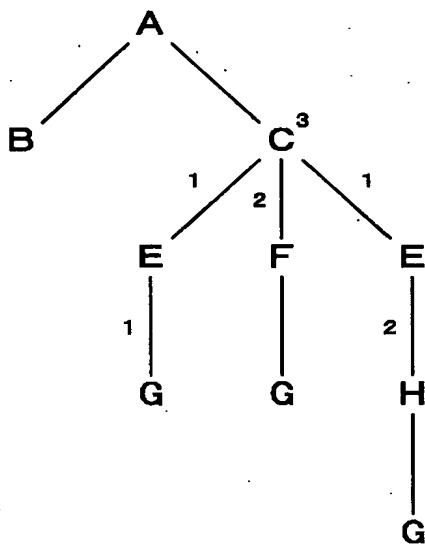
【図 7】



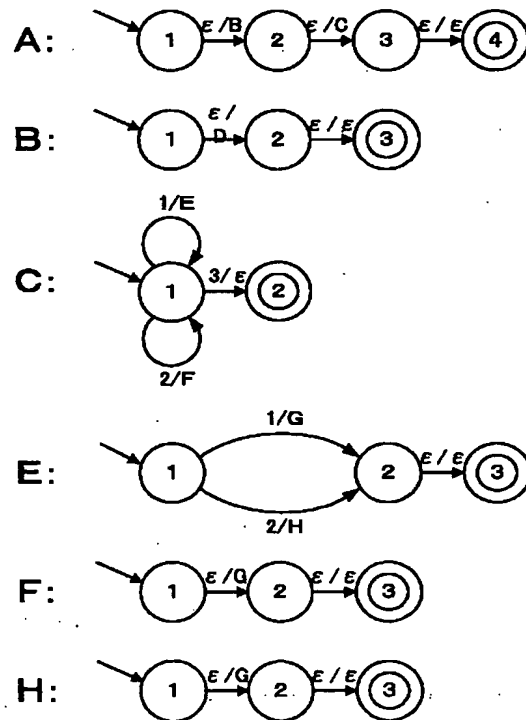
【図 9】



【図 10】



【図 11】

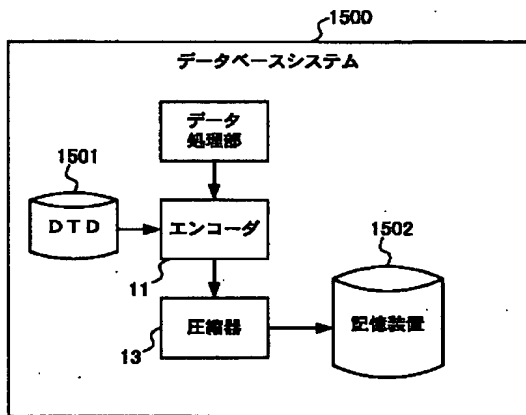


【図12】

DTD

変換前	変換後
<ELEMENT department (employee)>	<ELEMENT department (employee)>
<ELEMENT employee (name, email)>	<ELEMENT employee (serialNo, manager?, name, email)>
<!ATTLIST employee serialNo CDATA #REQUIRED>	<ELEMENT serialNo (#PCDATA)>
<!ATTLIST employee manager CDATA #IMPLIED>	<ELEMENT manager (#PCDATA)>
<ELEMENT name (#PCDATA)>	<ELEMENT name (#PCDATA)>
<ELEMENT email (#PCDATA)>	<ELEMENT email (#PCDATA)>

【図15】



フロントページの続き

(72)発明者 丸山 宏
 神奈川県大和市下鶴間1623番地14 日本ア
 イ・ビー・エム株式会社 東京基礎研究所
 内

(72)発明者 田村 健人
 神奈川県大和市下鶴間1623番地14 日本ア
 イ・ビー・エム株式会社 東京基礎研究所
 内

(72)発明者 浦本 直彦
 神奈川県大和市下鶴間1623番地14 日本ア
 イ・ビー・エム株式会社 東京基礎研究所
 内

Fターム(参考) 5J064 AA00 BA11 BA14 BC02 BD02
 BD03